

United States Patent Application

Title

METHOD AND APPARATUS FOR SEARCHING AND
DISPLAYING STRUCTURED DOCUMENT

Inventors

Takuya OKAMOTO,

Toru TAKAHASHI,

Yuki AOYAMA,

Noriyuki YAMASAKI,

Eiko MURATA.

RECEIVED "SA" 11/19/80

BACKGROUND OF THE INVENTION

The present invention relates to a technology of searching and displaying a structured document produced in the standard generalized markup language (SGML) or the hypertext markup language (HTML), or more in particular to a method and an apparatus for searching and displaying the result of searching a structured document in highlight.

With the extension of ownership of word processors and the like, the document information produced are going electronic more and more. These electronic documents have individual unique formats depending on the equipment or the software which has produced the documents and cannot be used with another equipment or software. The need has arisen, therefore, for some conversion means.

Various structured documents have been proposed as a common format for exchanging such documents. These structured documents can define the hierarchical structure including chapters, sections and paragraphs constituting a basic structure of documents and also can contain layout information.

A descriptive language for the structured documents for which standardization is under way is a standard generalized markup language (SGML). SGML uses

108250544950
a method of expressing a document element by embedding a
specific character string called a tag in the text as
element information of a structured document. According
5 elements indicated by tags can be defined by a document
type definition (DTD).

The above-mentioned SGML and DTD are described
in detail in "Practical, SGML" (edited and translated by
the SGML Gathering, Working Group for Practical Applica-
10 tion, April 20, 1992, published by Japan Standards
Association).

Assume that these structured documents are
registered in the data base of a search system and
searched by specifying an element name. In the case
15 where the DTD varies from one document to another to be
registered, a processing method is to analyze the
elements of each document, determine which portion of
the document corresponds to a specified element name,
and acquire and retrieve the character string to be
20 searched.

This method, however, consumes considerable
time for processing. Also, in a method using a table
listing a portion of each document corresponding to each
element name, it is necessary that all the element names
25 appearing in each document are managed collectively and
to register all corresponding portions of each document
for each element name. This requires a management table
of enormous size.

Further, all documents in registration with different DTDs do not necessarily have the same element to be searched. Also, in the case where different names of the same meaning such as "abstract" and "gist" are
5 attached to elements, all the different element names have to be specified for search. In actual practice, therefore, a structured document cannot be searched easily.

For the search of a structured document,
10 therefore, it is necessary to register only the documents generated according to the same document type definition. In this way, element names specified in advance are used to manage corresponding portions of each document.

15 At the time of search, an element name to be searched and a query are specified. If a character string meeting the query is contained in the portion of each document corresponding to the specified element, the query is judged as matching.

20 An explanation will be given of conventional techniques having the function of displaying the contents of a document as the result of searching a structured document.

A first conventional technique that can be
25 cited is JP-A-8-339369 entitled "Document display apparatus and document display method".

This conventional technique discloses a method of converting into a layout for element analysis and

0954445-09304
T08360-5449560

element display and displaying the contents of a specified element of a SGML document. It is possible to display a structured document by element using this technique. Further, this conventional technique provides means for highlighted display (an intensified display with the color, style or size of a character changed or a character underlined) of a specified element.

10 The means for highlighted display disclosed in this conventional technique, however, is for controlling a display method for each element, and specifies whether a particular element is displayed or not displayed and whether it is displayed in highlight or not. This conventional technique, therefore, fails to disclose a method of realizing highlighted display of a matching query term which is required for displaying the result of searching a structured document.

15 A second conventional technique disclosed in JP-A-8-212230 entitled "Method of document search and document searching apparatus" is a method for highlighted display of the result of searching a document other than a structured document.

20 This conventional technique, however, only acquires a matching strings position of a text for display and adds highlight information, but has no function of adding the highlight information to a document obtained as a result of searching a structured document.

5 A mere combination of these two conventional techniques cannot realize the function of adding the highlight information to a matching query term in a document output as the result of searching a structured document.

Specifically, highlight display of a structured document requires means for producing a DTD with element information for highlight added to the DTD used for producing a document to be displayed.

10 A method of altering the document type definition for adding highlight information to a structured document is disclosed in JP-A-8-159202 entitled "Method and apparatus for plate management of structured documents" constituting a third conventional technique,
15 in which a DTD is produced by adding a new element to the original DTD.

The use of this conventional technique makes it possible to produce a document type definition with the highlight information added thereto.

20 It is seen that the first and second conventional techniques permit a structured document to be displayed with the elements thereof clearly known on the one hand and permit a highlighted display of a matching strings position of a document not structured on the
25 other.

Further, the use of the third conventional technique makes it possible to specify a document type definition with highlight information added for each

0956443-095644

element.

By combining these techniques, it is possible to output a structured document with highlight information added to the result of searching a specified
5 element thereof and thereby to realize a highlighted display of the structured document.

In recent years, the internet has explosively spread as a method of acquiring the latest information. Also, the function of searching information on a web has
10 been improved as a means for quickly acquiring information required by the user from a great amount of information available on the internet.

The hypertext markup language (HTML) is for describing the contents of a document and expressing
15 information for linking to other resources and a document format on WWW (World Wide Web). HTML is regarded as a SGML described in accordance with a specified DTD. A means for producing and processing a HTML document is a HTML editor. A HTML browser, on the
20 other hand, analyzes and displays the HTML document thus produced.

There is a type of HTML browser which is supplied with a character string (hereinafter referred to as "the query term") and which has such functions as
25 searching a HTML document on display and displaying a matching strings position intensively by reverse video or the like.

0964475-092804

5

10

given structured document can be displayed individually in highlight by combining the above-mentioned conventional techniques.

15

and displaying the original structured document with

highlight information added thereto, the mere use of the above-mentioned method of the conventional techniques for the normalized structured document is not sufficient. In other words, since only a portion of the element information of the original document remains available at the time of search, the conventional method cannot realize highlighted display of the original structured document matching a query term simply by adding the highlight information to the element information.

An object of the present invention is to realize conversion from a document to be searched to the highlight position information of the original document in order to add the highlight information to the original document based on the result of searching a normalized document.

Another object of the invention is to realize a method and an apparatus wherein in the case where a matching query term after normalization covers a plurality of elements of the original document, highlight information is added to the matching strings position for each element to achieve highlighted display.

Still another object of the invention is to provide a method and an apparatus wherein in order to display in highlight the entire element including a matching query term or to display in highlight the entire area including two query terms satisfying the

proximity condition of the occurrence position, or in order to execute other similar processes, hierarchical highlight information is added for highlighted display according to different highlight display formats.

5 Yet another object of the invention is to provide a method and an apparatus wherein in the case where only a subelement of a structured document is extracted and displayed, the contents of such a subelement are displayed in highlight with highlight
10 information added thereto.

Some HTML documents are produced based on a plurality of DTDs by unique expansion dependent on the browser and it is difficult to determine a DTD on the basis of which a HTML is written. Further, there are
15 many HTML documents not correctly written according to the SGML grammar. It is therefore difficult to analyze the structure of the HTML document by the same method as the SGML document.

Other problems include:

20 (1) For a plain text document, a HTML document is produced with a highlight tag inserted before and after a matching strings position after search, so that a matching character string can be displayed intensively on a HTML browser. In the case where a character string
25 in a tag coincides with the query term, however, insertion of a highlight tag before or after the matching strings position would alter the contents of the tag in the original HTML and thus poses the problem that

0964475-092804
FOUO-544950

correct display is impossible.

(2) A tag for expressing a layout may be inserted amid a character string displayed continuously on the HTML browser. Correct search of the HTML document is impossible unless the tag is removed beforehand. Assume, for example, that the statement "This month's feature article" written in the HTML document, and that the query term is "feature article". In the HTML document, the tag "" for displaying a character in enlarged form is written between "feature" and "article". Thus, correct search is impossible unless the tag is skipped.

In order to solve the above-mentioned problems, according to a first aspect of the present invention, there is provided a method of searching and displaying a structured document for an information processing system including a processor, a memory unit, a file unit and an input/output unit, wherein the processor comprises: means for analyzing a structured document input thereto thereby to generate an analyzed structured document and storing the analyzed structured document in the file unit; means for searching the document search indexes stored in the file unit according to a query input thereto, determining whether or not there is content information meeting the query, acquiring an analyzed structured document having the content information considered to meet the query, and acquiring the information on the position of the document meeting

096446-09804
TOP SECRET 574950

the query; means for producing a document type definition (DTD) for highlighted display of the position of the document meeting the query; and means for producing a structured document for display with information on the position of the document meeting the query and information for highlighted display added to a structured document based on the document type definition for display.

According to a second aspect of the present invention, there is provided a method of searching and displaying a structured document for an information processing system including a processor, a memory unit, a file unit and an input/output unit, wherein the processor comprises: means for analyzing a structured document input thereto thereby to generate an analyzed structured document and storing the analyzed structured document in the file unit; means for generating a normalized structured document for document search with predetermined no-search element information removed from the input structured document, generating the restoring information for restoring the removed element information and storing the normalized structured document in the file unit; means for searching the normalized structured document stored in the file unit according to a query input thereto, determining whether or not there is a normalized structured document meeting the query, acquiring a normalized structured document considered to meet the query, and acquiring the information on the

00964435-092301

position of the document meeting the query; means for
producing a document type definition for highlighted
display of the position of the document meeting the
query; and means for restoring a structured document
5 having the removed element information based on the
restoration information from the normalized structured
document acquired by search, and producing a structured
document for display with information on the position of
the document meeting the query and information for
10 highlighted display added to the restored structured
document based on the document type definition for
display.

According to a third aspect of the present
invention, there is provided a method of searching and
15 displaying a structured document for an information
processing system including a processor, a memory unit,
a file unit and an input/output unit, wherein the
processor comprises: means for analyzing a structured
document input thereto thereby to generate an analyzed
20 structured document and storing the analyzed structured
document in the file unit; means for acquiring content
information in each element from the analyzed structured
document thereby to generate a document search index and
storing the index in the file unit; means for searching
25 the document search index stored in the file unit
according to a query input thereto, determining whether
or not there is content information meeting the query,
acquiring an analyzed structured document having the

0096445-002801

content information considered to meet the query, and
acquiring the information on the position of the
document meeting the query; means for acquiring an input
subelement to be displayed; means for producing a
5 document type definition for highlighted display of the
position meeting the query in the subelement to be
displayed; and means for producing a structured document
for subelement display with information on the position
of the document meeting the query and information for
10 highlighted display added to a structured document based
on the document type definition for subelement display.

According to a fourth aspect of the present
invention, there is provided an apparatus for searching
and displaying a structured document for an information
15 processing system including a processor, a memory unit,
a file unit and an input/output unit, wherein the
processor comprises: means for analyzing a structured
document input thereto thereby to generate an analyzed
structured document and storing the analyzed structured
20 document in the file unit; means for generating a
normalized structured document for document search with
predetermined no-search element information removed from
the input structured document, and storing the norma-
lized structured document in the file unit; means for
25 generating the restoring information for restoring the
removed element information and storing the restoring
information in the file unit; means for searching the
normalized structured document stored in the file unit

09964475-092304

according to a query input thereto, determining whether or not there is a normalized structured document meeting the query, acquiring a normalized structured document considered to meet the query, and acquiring the information on the position of the document meeting the query; means for producing a document type definition for highlighted display of the position of the document meeting the query; means for restoring a structured document having the removed element information based on the restoration information from the normalized structured document acquired by the search; and means for producing a structured document for display with information on the position of the document meeting the query and information for highlighted display added to the restored structured document based on the document type definition for display.

According to a fifth aspect of the present invention, there is provided a method of searching and displaying a structured document for an information processing system including a processor, a memory unit, a file unit and an input/output unit, wherein the processor comprises: means for storing a structured document, as a plain text with a tag, conforming with a specific document type definition input thereto; means for searching the plain text stored in the file unit according to an input query, determining whether or not there is a position meeting the query, acquiring the document having a position meeting the query as a plain

09964475-092801

test, and acquiring the information on a position of the document meeting the query; and means for producing a structured document for display with information added thereto for highlighted display of the position of the document meeting the query based on the specific document type definition as a document type definition for display.

According to a sixth aspect of the present invention, there is provided a method of searching and displaying a structured document for an information processing system including a processor, a memory unit, a file unit and an input/output unit, wherein the processor comprises: means for storing a structured document, as a plain text with a tag, conforming with a specific document type definition input thereto; means for searching the plain text stored in the file unit according to an input query, determining whether or not there is a position meeting the query, acquiring the document having a position meeting the query as a plain test, and acquiring the information on a position of the document meeting the query; means for determining whether or not a position meeting a query exists in the attribute information of a tag indicating a document element in a structured document; and means for adding, in the case where a position meeting the query exists in the attribute information of a tag, a character string including a character string of the position meeting the query to the content of the structured document, and

09964475-092804

producing a structured document for display with information added thereto for highlighted display of the position meeting the query in the character string based on the specific document type definition.

5 According to a seventh aspect of the present invention, there is provided a method of searching and displaying a structured document for an information processing system including a processor, a memory unit, a file unit and an input/output unit, wherein the

10 processor comprises: means for storing in the file unit a structured document based on a specified input document type definition as a plain text with a tag, and means for removing a character string constituting a predetermined specific tag from the character string to

15 be searched, coupling and searching the character strings before and after the character string constituting the specific tag to each other and searched, and adding information for highlighted display of the position meeting the query obtained by the search based

20 on the specific document type definition to the position meeting the query thereby to produce a structured document for display.

 According to an eighth aspect of the present invention, there is provided a method of searching and

25 displaying a structured document for an information processing system including a processor, a memory unit, a file unit and an input/output unit, wherein the processor comprises: means for storing in the file unit

09964475-092301

09964475-092804
TOP SECRET

a structured document based on a specified input document type definition as a plain text with a tag; means for searching a structured document stored in the file unit as a plain text according to an input query, and determining whether or not a position meeting the query is interposed between a specific tag indicating the start of a predetermined document element and a specific tag indicating the end of the document element; and means for adding, in the case where the position is so interposed, a character string including a character string of the position meeting the query to the content before the specific tag indicating the start of the document element or after the tag indicating the end of the document element, and producing a structured document for display with information added thereto for highlighted display of the position meeting the query in the character string based on the specific document type definition.

BRIEF DESCRIPTION OF THE DRAWINGS

20 Fig. 1 is a block diagram showing the process executed by an apparatus for searching and displaying a structured document according to first and second embodiments.

25 Fig. 2 is a flowchart for the process of searching and displaying a structured document.

 Fig. 3 is a diagram showing the registration of a structured document.

Fig. 4 is a flowchart for the process of registering a structured document.

Fig. 5 is a diagram showing a text for search.

Fig. 6 is a flowchart for the updating
5 process.

Fig. 7 is a flowchart for the process of extracting a specified element.

Fig. 8 is a diagram showing the information output as the result of analysis of a specified element.

Fig. 9 is a flowchart for the document display
10 process.

Fig. 10 is a diagram showing an example of a structured document and an example of the highlighting process.

Fig. 11 is a flowchart for the process of
15 producing a DTD for document display.

Fig. 12 is a diagram showing the process of normalization for searching a structured document.

Fig. 13 is a diagram showing the contents
20 stored as a result of the normalization process.

Fig. 14 is a diagram showing the process of conversion of the matching strings position information after normalization.

Fig. 15 is a flowchart showing the process of
25 conversion of the matching strings position information after normalization.

Fig. 16 is a flowchart for the process of adding highlight information.

0964475-09604
T08260-544950

Fig. 17 is a diagram showing the matching strings position information according to the second embodiment.

Fig. 18 is a diagram showing a definition of a highlighting method for each matching strings position information according to the second embodiment.

Fig. 19 is a diagram showing the conversion to the DTD for highlighted display according to the second embodiment.

Fig. 20 is a flowchart for the highlighting process according to the second embodiment.

Fig. 21 is a diagram showing an example of the SGML document with highlight information added thereto according to the second embodiment.

Fig. 22 is a diagram showing an example of highlighted display.

Fig. 23 is a block diagram schematically showing an apparatus for searching and displaying a structured document according to a third embodiment.

Fig. 24 is a flowchart for the process according to the third embodiment.

Fig. 25 is a diagram showing the process of conversion to the DTD for subelement display.

Fig. 26 is a flowchart for the process of producing the DTD for subelement display.

Fig. 27 is a diagram showing a system configuration according to a fourth embodiment.

Fig. 28 is a flowchart for a data controller.

099644E-099644

Fig. 29 is a flowchart for the process of character search and production of highlight position information according to the fourth embodiment.

5 Fig. 30 shows a configuration of a highlight position information storage area.

Fig. 31 shows a configuration of a highlight number storage area.

Fig. 32 shows a configuration of a highlight tag character storage area.

10 Fig. 33 is a flowchart for the process of producing a HTML document with highlight tag according to the fourth embodiment.

Fig. 34 shows an example of highlight insertion.

15 Fig. 35 shows an example of highlight inserted.

Fig. 36 shows a system configuration according to a fifth embodiment.

20 Fig. 37 is a flowchart for the search process and the process of producing the highlight position information according to the fifth embodiment.

Fig. 38 is a flowchart for the search internal to a tag and the search external to a tag according to the fifth embodiment.

25 Fig. 39 is a flowchart for the search external to tag according to the fifth embodiment.

Fig. 40 is a flowchart for the process of producing a HTML document with highlight tag according

2025051409550

Fig. 41 is a flowchart for the process of highlight tag insertion according to the fifth embodiment.

5 Fig. 42 is a flowchart for external highlight
tag insertion according to the fifth embodiment.

Fig. 43 shows an example of query according to a sixth embodiment.

Fig. 44 shows an example of a matching strings
10 position information according to the sixth embodiment.

Fig. 45 is a diagram showing the process of converting to the DTD for highlighted display according to the sixth embodiment.

Fig. 46 is a diagram showing an example of the
15 SGML document for highlighted display according to the
sixth embodiment.

Fig. 47 is a diagram showing an example of highlighted display according to the sixth embodiment.

DESCRIPTION OF THE PREFERRED EMBODIMENTS

20 A block diagram of a first embodiment is
schematically shown in Fig. 1.

Reference numeral 101 designates a structured document search and display apparatus. With a structured document 102 stored in a registered data file 114 as an input, a document is registered thereby to generate an analyzed structured document (described later with reference to Fig. 3) and a search index

(described later with reference to Fig. 5).

The analyzed structured document is stored in a structured document data base (hereinafter referred to as the data base DB) 105, and the search index is stored
5 in a search index DB 106.

Then, a query 103, which is applied from an input/output unit 115, is analyzed and the search index is read out for executing a search process 108. As a query result, information 109 of a matching document
10 identifier and information 110 of a matching strings position are output.

In the display process, first, a specified analyzed structured document 111 is read out from the structured document DB 105 based on the matching
15 document identifier information 109 in a document read process 107. In the document display process 112, a structured document 113 for display with matching information embedded therein is generated from the analyzed structured document 111 based on the matching
20 strings position information 110. The structured document for display thus generated is displayed on the input/output unit 115.

Fig. 2 shows a flowchart of the process for search and display of a structured document.

25 First, a structured document is registered (201). The registration process will be described later with reference to the flowchart of Fig. 4.

09064475-092804
F08260-52449660

Then, a structured document is searched using a specified query (202). The search process will be described in detail later with reference to the flowchart of Fig. 6.

5 The query result includes the number of matching documents, the number for identifying a matching document, and the matching strings position of the query term for each document. The matching strings position information output include an element ID (element
10 identifier) for identifying the element containing the matching query term, the matching start position in the element and the information on the text length.

 In the case where the number of matching documents is 1 or more in the process of structured
15 document search (203), the contents of the matching document are read (204), the matching strings position information of the document read is acquired (205) and the highlighted display is realized (206) in that order. The display process will be described in detail later
20 with reference to Fig. 9.

 In the presence of another matching document, the steps 204 to 206 are repeated.

 Upon completion of the display process, the presence or absence of the next search process is
25 checked (208). In the absence of a query, the process is terminated, while in the presence of a query, the process is returned to step 202 for repeating the search and display of a structured document.

Fig. 3 is a diagram schematically showing the process for registration of a structured document.

First, the elements of a SGML document are analyzed, and a structure tree 302 thereof is generated.

5 The contents of each item of the structure tree thus generated are output as data 303 of table format and registered as an analyzed structured document. In the process, CDATA designates character string data.

Fig. 4 is a flowchart for the process of registration of the structured document.

10

First, the structured document is analyzed (401). The structured document thus analyzed is registered as an analyzed structured document (402). For analysis of a structured document, a SGML purser can

15 be used for analyzing the SGML document using DTD.

Then, the analyzed structured document is normalized (403) for removing the elements not required for search.

The normalization process will be described later with reference to Fig. 12. The normalized structured document is registered in a document data base (404).

20

Further, from the analyzed structured document registered in the data base, the element information and the information on the text in the element and are

25 retrieved (405) as search indexes required for searching the structured document. The search indexes thus obtained are registered in a search index data base

(106). The search indexes thus registered lack the element information (tag) in the SGML document and are a stored text string indicating the element information and the contents of each element.

5 Fig. 5 shows an example storage of a search text including the search index and a normalized structured document. The above-mentioned process is repeated for the registered documents until the registered documents are depleted (407) when the process is terminated.

10 The contents in registration are used for full-text search of the registered documents.

 Fig. 5 shows an example of contents output as a text for search. This information including a table containing the element ID of the document structure and

15 a corresponding text string and character string information as shown in this diagram is registered as a text for search. At the time of search, a character string required for search is extracted based on the element ID.

20 Fig. 6 shows a flowchart for search of a structured document in step 202 of the structured document search and display process shown in Fig. 2.

 The query is given in such a form as "Specify an element to be searched: query expression".

25 Each element to be searched is defined by, for example, "<" and ">" like "<document.title>", and a containing element ("document" in the shown case) and a subelement ("title" in the shown case) are discriminated

by ".", thereby specifying an element to be searched in a hierarchy structure.

The query expression "and("search","document")", for example, indicates the condition under which both
5 "search" and "document" occur, and $C \leq 10$ ("search","document") indicates the condition under which "search" and "document" occur with ten or less characters interposed therebetween.

For search of a structured document, first, a counter of the number of documents is cleared (601), and
10 then the elements specified to be searched in the query are analyzed (602). At step 602, an element ID (element identifier) that can uniquely specify an element corresponding to an analyzed structured document is acquired from a character string specifying an element such as
15 <document.title>. The process for acquiring an element ID will be described later with reference to the flowchart of Fig. 7.

As the next process, a document (text for search) registered for search is read out, and the text
20 portion corresponding to the specified element ID acquired at step 602 is acquired (603).

A query expression including a query term, the AND of a plurality of query terms that occur, a proximity condition and the like logical conditions is
25 analyzed (604) based on the query, and the query term thus obtained is used to make full-text search of the text portion acquired at step 603. Thus it is determined whether the logical conditions of the query

Once the query is matched (606), such information as the document index, the ID of the element
5 containing the query term and the position where the
query term matches in the element are output (607).

Fig. 7 is a flowchart showing the process for analyzing the structure-specified contents in the query analysis of Fig. 6.

20 In the presence of a subelement (705), it is
further determined whether a still lower subelement, if
any, has a structure specified in similar fashion. In
the case where the answer is affirmative, step 706 for
outputting an element ID is repeated until the subele-
25 ments are depleted (707). Upon completion of the
process for all the elements, a list of the element IDs
to be searched is produced.

The number of the element IDs to be searched
(801) and the IDs in the number obtained for search
5 (802) are output.

First, the structured document to be searched is a normalized one in which the elements not required for search are removed. The elements matched by search or the matching strings position information are not necessarily coincident with those of the registered original structured document (See the structure tree 302 of Fig. 3 and the structure tree of Fig. 12).

Consequently, first, the DTD for display of
20 the document to be displayed is produced from the DTD of
the registered document (901). The process for produc-
ing the DTD for display will be described later with
reference to Fig. 11.

Further, the matching strings position
25 obtained for the normalized structured document is
converted into the information on the highlight position
and the elements for the registered original structured
document (902). The process for converting the matching

position information of the normalized document into the highlight position information of the original document will be described later with reference to Fig. 15.

Then, the information of the base document
5 element of the analyzed document used for display are read out, and steps 903 to 911 are repeated sequentially thereby to output a document for display.

First, the element information is read out (903), and then an element start tag is output (904).
10 Further, in the presence of at least a subelement in the element (905), the display process is recursively carried out (906) for the subelements (steps 903 to 911). After depletion of the subelements, the process proceeds to step 911 for outputting the tag indicating
15 the end of the element.

The subelement includes a character string.
Therefore, such a structured document as

```
<document>
<title>
20 structured document
</title>
<text>
<intensify> structured document <intensify> is
searched.....
25 </text>
</document>

has an element in terms of a character string (expressed
as CDATA for SGML) as a subelement of <title>. CDATA
```

has no subelement, and constitutes a character string information which has the contents of "structured document" in the shown case.

Similarly for <text>, the element <intensify>
5 and the character string having the content "is searched....." exist as subelements.

In the case where it is determined that no subelement exists at step 905, the element is a character string. Therefore, the matching position
10 information is compared with the contents of this element (908), and in the case where the element includes a matching position, the highlighting process is carried out (909). The highlighting process will be described later with reference to Fig. 16.

15 In the case of the character string contains no matching position, on the other hand, the content is directly output as a text (910). In the case where the output content is a character string, the start tag or the end tag is not output at steps 904, 911.

20 The highlighted display is realized for each element by the above-mentioned process. In the presence of any other elements to be further processed, the process from step 903 is repeated (912).

Fig. 10 shows a DTD for registration (1001),
25 an example (1002) of the SGML document (document instance) to be registered, a display DTD (1003) used for highlighted display and an example (1004) of the SGML document (document instance) converted for display.

character string is mixed with an element, the use of
#PCDATA is required.

As a content model, "RCDATA" instead of
"CDATA" may be specified: The difference between RCDATA
5 and CDATA is that in the case where a reference to
entity (described like "&xxxx". Used for replacement
with an exceptional character or the like) occurs in
the element, the character string that occurs is handled
as it is without conversion to entity (exceptional
10 character or the like). In the case where "RCDATA" is
specified, a character string that has been converted to
entity is handled as such.

For highlighted display, the document
structure is required to be altered to permit highlight
15 information to be added to the character string. The
element information for highlighted display is added to
all the character strings of each element like points of
alteration underlined at 1003, to which the element
information for highlighted display "<!ELEMENT highlight
20 ..(#PCDATA)>" is required to be further added.

The portion "CDATA" in the content model of
the original DTD is replaced by "(#PCDATA!highlight)*"
because CDATA indicates that only one character string
exists in the element thereof and cannot occur as a
25 repetitive element. Since a tag for highlight is added,
CDATA in the original element is altered to #PCDATA,
and then altered to "(#PCDATA!highlight)*" to permit
repetitive occurrence of highlight.

Fig. 11 is a flowchart showing the process for producing the DTD for highlighted display from the DTD for registration.

First, the DTD for registration is read (1101) and the content of the DTD is analyzed to acquire the ELEMENT items (1102). In the case where CDATA, RCDATA, #PCDATA or the like is specified in the content model of the ELEMENT item, the content models are all altered in such a manner that the element for highlight can be added (1103 to 1106).

For altering the content model, first, "CDATA", "RCDATA" are altered to "#PCDATA", after which "#PCDATA" is defined in such a form as "(#PCDATA!highlight)*" so that a character string surrounded by the highlight tag and a character string not so surrounded may appear repeatedly.

In the case where the original content model is described as "(#PCDATA!underline)*" in such a manner that a plurality of elements may occur repeatedly, the description "(#PCDATA!underline!highlight)*" is sufficient to indicate the occurrence of a highlight element.

Upon complete alteration of all the ELEMENT declarations (1107), "<!ELEMENT highlight..CDATA>" is added as a definition of the element for highlight (1108). The foregoing process generates a DTD for highlighted display shown by 1003 of Fig. 10.

Fig. 12 shows the process for normalizing the

structured document.

The structured document designated by 1001 of Fig. 10 can be expressed by a structure tree of 1201.

In the case where "underline" is specified as
5 an unrequired element, the first step of normalization
is to delete the element "underline" as shown in 1202
while the character string contained in the subelement
of the underline is incorporated as an element of the
"text" constituting an immediate containing element.

10 Further, the two character strings (CDATA)
existing as subelements of the "text" are coupled into a
single character string as shown in 1203.

Fig. 13 shows the original structured document
(1301) and the normalized structured document (1302)
15 whose contents are analyzed and converted into and
output as a table. Numeral 1303 designates a table
storing the element information, in which the elements
with the element IDs of 0 to 6 are the information on
the original elements. Numeral 0 is the base document
20 element, and the document structure can be determined by
tracing the information of subelements.

The elements with the element IDs (element
identifier) of 7 to 9 attached thereto are those altered
and added after normalization.

25 Numeral 7 designates the base document
element, and the normalized document structure can be
determined by tracing the subelements. The element
information of the elements ID1, ID2 including "title"

and the underlying elements not altered are left as they are.

Further, the correspondence between the elements ID7 to ID9 added by normalization and the original elements thereof is stored in a normalization
5 correspondence table of 1304.

Fig. 14 shows the result of converting the information on the matching strings positions searched for the normalized structured document into the position
10 information for the original structured document.

The information 1401 on the matching strings positions obtained from the normalized elements is converted into the position information 1402 for the original structured document using the information in
15 the normalization correspondence table 1304 in Fig. 13.

In the shown example, the matching strings position of the element ID9 after normalization is divided into the elements ID5 and ID6 for the original document, and therefore is altered to the position
20 information to be highlighted in the two elements.

Fig. 15 is a flowchart showing the process for converting the matching strings position information of the normalized structured document in step 902 of Fig. 9 into the matching strings position information for the
25 original structured document.

First, the matching strings position information of the normalized structured document are sequentially read (1501), and it is determined whether

In the case where the element IDs exist from before normalization, there is no alteration, and therefore the matching strings position information before normalization is output as it is (1503).

15 Once the matching strings position is obtained
for an element in the original structured document, it
is output as a matching strings position in the original
structured document (1505).

Fig. 16 is a flowchart for the highlighting process of step 909 in Fig. 9. First, the character string from the document head to the highlight start is output (1601). Then, the start tag of the element used for highlighted display is output (1602).

Further, the character string in highlight position is output (1603) and the end tag for the

element used for highlighted display is output (1604).

Upon completion of the entire highlighting process (1605), the remaining text is output thereby to end the highlighting process (1606).

5 Now, a second embodiment will be explained with reference to the process for altering the highlighted display method according to the matching condition and the process for executing a plurality of highlighting processes hierarchically. The block
10 diagram schematically showing the process is the same as Fig. 1.

Fig. 17 shows the matching strings position information (1701) used in the present embodiment.

15 The information added to the matching strings position information shown in Fig. 14 represents the area 1702 added for storing each condition matched.

20 Further, although only the position of the matching query term is output in Fig. 14, this embodiment makes it possible to specify an area including the query term by highlighting the whole element containing the query term in addition to the matching query term according to the query.

25 These information on matching conditions are added at the time of searching the structured document. In the case under consideration, such indexes as the proximity condition used for the query and the frequency of occurrence of each query term are added. Alternatively, however, each condition can be weighted for

each query term in advance.

Fig. 18 is a table 1801 defining the correspondence between the matching condition and the highlighting method (form of highlighted display).

5 Highlighting methods 1803 corresponding to the matching conditions 1802-are described. The position matched according to each matching condition is displayed in highlight based on the contents of this table.

10 Further, hierarchy information 1804 is given. The larger the value of the hierarchy information for an item, the higher the level of highlighting the particular item, such as when highlighting the whole element.

Fig. 19 shows the process for producing a DTD
15 for display to realize the above-mentioned highlighting process. Based on the original DTD 1901 used for registration, the DTD 1902 for highlighted display is generated, in which the definition in the high-level highlight element is altered or added to make it possible to specify or omit a low-level highlight element
20 hierarchically.

In producing the DTD, a plurality of highlight information in the above-mentioned process shown in Fig. 11 are all added (1903) when adding the highlight
25 information at step 1106. Further, when adding the ELEMENT declaration for highlight at step 1108, the low-level highlight elements and character strings are incorporated as a content model constituting subelements

of each highlight element based on the hierarchical information 1804 of Fig. 18.

In the absence of a low-level highlight element, only a character string occurs (1904) as a content model.

Fig. 20 is a flowchart for the highlighting process according to the second embodiment.

First, the highlight information are sorted with the order of the starting position as a first key and the order from upper to lower level in hierarchy information as a second key (2001). Then, the text up to the highlight start is output (2002), and a highlight start tag is output (2003).

Further, if the next highlight is started before the end of a highlight position, it indicates the presence of a low-level element information (2004). Thus, the text up to that position is output (2005), after which the highlighting process is carried out for the low-level highlight element (2006). The highlighting process for the low-level subelement is the same as the process of steps 2003 to 2009.

If there is any lower-level highlight element (2007) at the end of the process for a low-level highlight element, the process is returned to step 2005 for outputting the text up to the next highlight element so that the lower-level highlight element is processed.

In the absence of any lower-level highlight element, the text up to the last low-level element is

output (2008) and a highlight end tag is output (2009).

In the case where there remains any information to be highlighted, the process is returned to step 2002 and repeated. Once the information to be highlighted is depleted (2010), on the other hand, the remaining text is output to end the process (2011).

Fig. 21 shows an example of the SGML document generated by the above-mentioned process.

Fig. 22 shows an example display of a text of the SGML document of Fig. 21. An overlapped highlight position is processed by repeating a highlighting method a plurality of times. An explanation will be given of the process for displaying in highlight by cutting out only a subelement of the structured document according to the third embodiment.

Fig. 23 is a block diagram schematically showing such a process according to the present embodiment.

The difference from Fig. 1 is that an element 2301 to be displayed is specified and that a subelement display process (2302) is executed instead of the document display process (112) based on the content of the specified element to be displayed.

Fig. 24 is a flowchart showing the sequence of the process for extracting and displaying a subelement.

First, a DTD for subelement display is generated (2401). The process for generating the DTD for subelement display will be described later with

Further, the matching strings position information obtained for the normalized structured document is converted into the element ID and the matching strings position information for the original document registered (2402). The process for converting the position information of the normalized document into that of the original document can use the method described above with reference to Fig. 16.

First, the element information to be displayed
15 is read (2403). It is determined whether or not the
particular element is to be displayed by use of the
method described above with reference to Fig. 7.

25 In the case where it is determined at step
2405 that there is no subelement, the element is that of
a character string. Therefore, the content of this
element is compared with the matching strings position

information (2408), and if the element contains a matching strings position, the highlighting process is carried out (2409). The highlighting process uses the method described above with reference to Fig. 15.

5 In the case where the character string contains no highlight position, the content is directly output as a text (2410). In the case where the output content is a character string, neither the start tag nor the end tag is output in steps 2404, 2411.

10 The highlighted display is realized for each element by the above-mentioned process. In the presence of any other element to be processed, the process from step 2403 is repeated (2412).

15 Fig. 25 shows the contents of the DTD to be produced for subelement display.

 In the subelement output, an element defined to always occur in the original DTD (2501) may not be output. Also, a containing element is not necessarily output.

20 As a result, it is necessary to change the process in such a manner that the occurrence of a start or end tag is not essential for a containing element and subelements may not necessarily occur. The DTD for subelement display thus produced is shown in 2502.

25 The SGML document produced using this DTD is shown in 2503. In this example, only the title is extracted.

 Fig. 26 is a flowchart showing the sequence of

generating a DTD for subelement display. First, the DTD for registration is acquired (2601).

Then, the ELEMENT items in the DTD are retrieved (2602). In the case where the content model
5 includes CDATA, RCDATA or #PCDATA, the highlight information is added (2603 to 2606).

The highlight information is added in the same manner as the process of steps 1103 to 1106 of Fig. 11.

Then, the occurrence indicators (*, +, ?, nil)
10 in the content model are checked. The indicator, if "+" (2607), is altered to "*" (2608). In the absence of an occurrence indicator (2609), "?" is added (2610).

Upon complete processing for all the ELEMENT declarations (2611), the ELEMENT declaration for the
15 highlight element is added (2612). Further, if the occurrence of the tag of an element having a subelement is "essential" (ü]), the indicator is altered to "unrequired" (0).

Now, a fourth embodiment of the invention will
20 be described with reference to the accompanying drawings.

Fig. 27 is a diagram showing a system configuration of this embodiment.

A WWW (world wide web) search system 2700 is
25 connected to a client 2701 through a network 2702. The client 2701 is a PC (personal computer), a WS (work station) or the like, and a query term is input on the query term setting screen on the web browser 2703

operating at the client 2701. The WWW search system 2700 makes search using this query term, and outputs the result of search to the web browser 2703.

The WWW search system 2700 includes a HTTP
5 server 2704 for receiving the query term from the client 2701, a data controller 2705 for conducting a searching operation and inserting a highlight tag, and a memory 2706 for storing the positional information of the highlight tag. The WWW search system 2700 is connected
10 to a magnetic disk drive 2707 for storing the HTML document to be searched.

The data controller 2705 searches the HTML document in the magnetic disk 2707 using the query term received from the HTTP server 2704, and inserts the
15 highlight tag at the matching strings position of the HTML document matched with the query term.

The memory 2706 includes a highlight number storage area 2708 for storing the number of matchings for each document, a highlight position information
20 storage area 2709 for storing the query result position information, a highlight tag character storage area 2710 for storing the contents of the highlight tag inserted, a HTML document temporary storage area 2711 for storing the HTML document with the highlight tag inserted
25 therein, and a query term storage area 2712 for temporarily storing the query term input by the client 2701 and acquired by the HTTP server 2704 of the WWW search system 2700.

09964475-092804
T09250-5249560

The HTML document with the highlight tag inserted therein by the WWW search system 2700 is displayed on the web browser 2703 of the client 2701 through the network 2702 from the HTTP server 2704.

5 Now, the process of the data processor 2705 will be specifically explained with reference to Fig. 28.

10 The query term set by the client 2701 is acquired and used for the search process, the matching strings position is detected, and highlight position information 2709 is produced. The highlight tag is embedded at the matching strings position of the HTML document matched with the query term, and displayed on the web browser 2703 of the client 2701.

15 Step 2800:

 The query term set by the client 2701 is acquired by the WWW search system 2700 using the HTTP server 2704. The query term thus acquired is stored in the query term storage area 2712 of the memory 2706.

20 Step 2801:

 The HTML document stored in the magnetic disk drive 2707 is subjected to full-text search using the query term stored in the query term storage area 2712 at step 2800. In the case of matching, the matching
25 strings position and the number of matchings in the HTML document are acquired, and the particular information are stored in the highlight position information storage area 2709 and the highlight number storage area 2708.

09964475-092881
T08260-52449550

This process will be described in detail with reference to Fig. 29.

Step 2802:

5 The highlight tag stored in the highlight tag
character storage area 2710 is inserted in the matching
strings position and stored in the HTML document
temporary storage area 2711 based on the information
stored in the highlight position information storage
area 2709 produced at step 2801. This process will be
10 described in detail with reference to Fig. 33.

Step 2803:

15 The HTML document for highlight stored in the
HTML document temporary storage area 2711 produced at
step 2802 is displayed on the web browser 2703 of the
client 2701 using the HTTP server 2704.

20 The process of steps 2800 to 2803 is repeated,
so that the HTML document stored in the magnetic disk
2707 is searched using the query input by the client
2701. Thus, a plurality of matching strings positions
for the document matched with the query can be displayed
in highlight.

Now, an explanation will be given of the
process for producing the highlight position information
of step 2801 in Fig. 28 with reference to Fig. 29.

25 Step 2900:

 The HTML document stored in the magnetic disk
2707 is read out. The HTML document 3400 of Fig. 34 is
an example thus read out.

0995443-092804
T08260"5449550

This HTML document 3400 is displayed on the screen of the web browser such as designated by 3401.

Step 2901:

5 The highlight position information storage area 2709 for storing α cases of highlight position information is secured, where α is an arbitrary positive integer. Also, a highlight number storage area 2708 for storing the number of highlights is secured.

10 The data formats of the highlight position information storage area 2709 and the highlight number storage area 2708 are shown in Figs. 30 and 31, respectively.

15 The highlight position information storage area 2709, as shown in Fig. 30, is configured of a HTML document identifier 3000, a highlight position number 3001 as counted from the head, a number 3002 of highlight bytes and a highlight tag number 3003.

20 The HTML document identifier 3000 is the number of the HTML document read at step 2900. The serial number or the like attached to the HTML document at the time of storage is stored as the HTML document identifier 3000.

25 The highlight position number 3001 indicates the matching string position in terms of the number of bytes as counted from the head of the HTML document read out at step 2900 and matched with the query term acquired at step 2800.

 The number 3002 of highlight bytes is stored

202505040950

in the form of the length highlighted in terms of the number of bytes. In other words, the length of the character string of the query term is stored.

5 The highlight tag number 3003 can discriminate the highlight tag for each of a plurality of query terms which may be used for highlighted display. The highlight tag is discriminated based on the information stored in this field. In other words, the data for discriminating the type of the tag used for highlighted
10 display is stored in this field.

Step 2902:

This step initializes the count *i_cnt* stored in the highlight position information storage area 2709.

Step 2903:

15 This step checks whether or not the query term read at step 2800 is coincident with the HTML document read at step 2900. In the presence of a matching point, the process proceeds to step 2904. In the absence of a matching point, on the other hand, the process proceeds
20 to step 2908. Step 2904:

This step checks whether or not the number stored in the highlight position information storage area 2709 secured at step 2901 or 2905 is larger than "*i_cnt*" indicating the number of highlights stored. In
25 the case where there still exists an area for storing data, the process proceeds to step 2906. In the absence of such an area, on the other hand, the process proceeds to step 2905.

09644-0304
100250-44950

Step 2905:

The highlight position information storage area 2709 is enlarged by a predetermined value and secured again, followed by proceeding to step 2906.

5 Step 2906:

The HTML document identifier 3000, the position 3001 as counted from the head of the HTML document, the number of highlight characters 3002 and the highlight tag number 3003 are stored at the (i_cnt)th
10 position of the highlight position information storage area 2709 secured at step 2902 or 2905. Since the count i_cnt is initialized to 0, the data are stored at the 0th position in the case where i_cnt is 0.

In the case where a plurality of highlight
15 information are stored in a single HTML document, i_cnt is updated and therefore the highlight information are stored at the position indicated by i_cnt.

Assume that the HTML document 3400 read at step 2900 is a HTML document identifier "001", and that
20 the query term extracted at step 2800 is a "feature".

When the query term "feature" is searched for in the HTML document 3400, the characters "feature" can be found at the 122nd byte (3403) as counted from the head of the HTML document 3400.

25 In this case, "001" (3404) is stored as the HTML document identifier 3000, "122" (3405) is stored at the position 3001 as counted from the head of the HTML document, and the number "4" of bytes (3406) for the

000000-000000

"feature" is stored as the number of highlight characters 3002. Finally, the number indicating the tag for intensifying the result of search is stored as the highlight tag number 3003. Such a number is "1" (3407)
5 in the case under consideration.

Fig. 32 shows a configuration of the highlight tag numbers and corresponding highlight tags actually stored. A structure 3200 for highlight tag insertion stored in the highlight tag character storage area 2710
10 is shown in (1) of Fig. 32.

The structure 3200 for highlight tag insertion is comprised of a tag number 1 (3202) for storing a serial number, a start tag 1 (3203) for storing the highlight start tag name, an end tag 1 (3204) for
15 storing the highlight end tag name, and a highlight tag number (3201) for storing the number of tags. There exist tag numbers, start tags and end tags in the number corresponding to the number of highlight tags stored in the highlight tag number field.

20 An actual example of the highlight tag character storage area is described in (2).

The description that follows concerns the case in which three types of highlight tags are stored. Therefore, "3" (3205) is stored in the area for storing
25 the number of highlight tags. A tag "<FONT COLOR =
"RED">" (3207) indicating red is stored in the start tag with the tag number "0" (3206), and "" (3208) as an end tag. In similar fashion, a tag "<BLINK>"

indicating a flicker is stored in the tag number "1" (3209), and "<H1>" for displaying the characters in enlarged size is stored in the tag number "2" (3210).

5 The highlight tag character storage area 2710 is produced before the highlight position information storage area 2709. The highlight character storage area 2710 can also be produced using the user interface.

10 In searching for a soundex or a synonym, therefore, provision of a plurality of highlight tags makes possible different highlighted displays for different queries by attaching the tag number "1" for the character searched in soundex and the tag number "2" for the character searched in synonym, for example.

15 In the case where "<BLINK>" is used as the highlight tag, "1" is stored as the highlight tag number 3407 in the highlight position information storage area 3402.

Step 2907:

20 Since data are stored in the highlight position storage area 2709 at step 2906, 1 is added to i_cnt and the process returns to step 2903.

Step 2908:

25 The number of highlights in the HTML document acquired at step 2900 is acquired and stored in the highlight number storage area 2708. The contents of the structure of the highlight number storage area 2708 will be explained with reference to Fig. 31.

Fig. 31 shows the contents of the structure of

09964475-092804
T08260"5449560

the highlight number storage area 2708. Numeral 3100 designates the document identifier of the HTML document read at step 2900. Numeral 3101 is the position where the acquired number of highlights is stored. In the
5 case under consideration, the document identifier "001" is stored as the document identifier 3100, and i_cnt is stored in the highlight number storage area 3101 thereby to end the process.

Now, the process for producing the HTML document with highlight tag will be explained with reference
10 to Fig. 33.

Step 3300:

This step checks whether or not it is necessary to insert a highlight tag in the HTML document
15 read at step 2900.

In the presence of any HTML document identifier 3000 stored in the highlight position information storage area 2709, the process proceeds to step 3301. In the absence thereof, on the other hand, all the texts
20 are output at step 3309 thereby to end the process.

Step 3301:

The process count i_cnt is initialized to 0.

Step 3302:

The HTML document temporary storage area 2711
25 is secured for storing the HTML document with a highlight tag inserted therein.

As the HTML document temporary storage area 2711, an area is secured in a size corresponding to the

5 The use of the above-mentioned process makes
it possible to search the HTML document based on the
query term set by the client 2701 and, for the document
coincident with the query term, to produce the contents
of the highlight number storage area 2708 for storing
0 the number of highlights and the highlight position
information storage area 2709 for storing the highlight
position.

This HTML document is displayed like 3503 with
the matching "feature" (3504) flickering.

25 Now, a fifth embodiment of the invention will
be explained with reference to Figs. 36 to 42.

Fig. 36 is a diagram showing a system configuration for a highlighted display method in which the

The process for search and production of
5 highlight position information is carried at step 2801
using the query term acquired at step 2800. The process
is specifically shown in the flowchart of Fig. 37.

The HTML document to be processed is read out
10 of a magnetic disk 2707.

The highlight position information storage area 2709 for storing the highlight position information and the highlight number storage area 2708 are secured in the memory 2706.

A highlight tag to be inserted before and after the matching string position is read out.

As seen from the specific example of applica-
tion shown in (2) of Fig. 32, the highlight tag is read
out of the highlight tag character storage area 2710.
In this case, the number of the highlight tag identi-
fiers is seen to be three from "3" (3205). The first
"0" (3206) has stored therein ""
(3207) and "" (3208). Thus, the start tag of the
highlight tag number 0 is "" and the
end tag "". In similar fashion, the start tag of
the highlight tag number 1 is "<BLINK>" and the end tag

The skip process is made possible by setting
5 the skip tag name storage area 3600 in advance of the
search process.

The number of characters found coincident with the leading character of the query term from the head of the HTML document at step 3704 is temporarily secured in the start position storage area 3601.

This step checks whether or not the character string of the query term is coincident with the characters written in the HTML document, and in case of coincidence, checks whether or not the point of coincidence exists inside or outside the HTML tag. Further, the position of the last character of the matching character string is secured by the number of characters as counted from the head of the HTML document. This process will be explained in detail with reference to Fig. 38.

25 This step checks for a matching as a result of
step 3706. In the case where a query term is existent
in the HTML document, the process proceeds to step 3708.
In the absence of a query term, on the other hand, the
process proceeds to step 3712.

Step 3708:

The highlight number storage area 3708 secured at step 3701 is compared with the number of stored highlights, and if the secured area is larger than the number of highlights, the process proceeds to step 3709. Otherwise, the process proceeds to step 3710.

Step 3709:

For lack of the area for storing data in the highlight position information storage area 2709, the area is set again and the process proceeds to step 3710.

Step 3710:

The number of characters to be highlighted and the information on the highlight position are stored in the highlight position information storage area 3600. Specifically, the document identifier of the HTML document read at step 3700 is stored as the HTML document identifier 3000 of the highlight position information storage area 3600 described with reference to Fig. 30, and the start position acquired at step 3705 is stored as the highlight-position-from-head information 3001. Also, the character string length of the query term is stored in the number of highlight bytes 3002, and the tag number read at step 3702 is stored in the highlight tag number 3003.

The highlight tag number 3003 has set therein "0" as a default value.

Step 3711:

In the case where there are a plurality of

0564479 092891
158560 544959

character strings matching with the query term, the process is executed to check for a point where the query term again matches with any of the characters following the first matching position in the HTML document. Thus, the sum of 1 and the number of characters from the head of the HTML document at the position where the last matching character secured at step 3706 is substituted into i_cnt. After updating the processing position, the process returns to step 3704.

10 Step 3712:

In the case where the character string from the start position stored in the start position storage area 3600 acquired at step 3705 fails to coincide with the query term, the process is executed to check for a point in the HTML document where the query term coincides again with any of the characters following the start position. The sum of 1 and the start position stored in the start position storage area 3600 is substituted into i_cnt. After updating the processing position, the process returns to step 3704.

The foregoing description concerns the process of searching including the checking inside and outside the tag and the process of producing the highlight position information.

25 Now, an explanation will be given of the process of searching inside and outside of the tag at step 3706 with reference to Fig. 38. In the process, it is checked whether or not the matching start position

acquired at step 3705 exists inside or outside the attribute of the tag indicating the document structure, and also it is checked whether or not the character string from the matching start position coincides with
5 the query term.

Step 3800:

This step checks whether or not the matching start position stored in the start position storage area 3600 at step 3706 is inside or outside the HTML tag.

10 The data are checked from the (i_cnt)th byte of the HTML document at the time point of step 3706 to the matching start position. The tag end character ">" corresponding to the tag start character "<" is checked thereby to check whether or not the matching start
15 position exists in the tag. In the case where there exists the tag start character "<" and the matching starting position is located before the tag end character ">", the starting position is assumed to exist in the tag, and the process proceeds to step 3801. In
20 the case where the matching starting position exists at the position not surrounded by the tag start character "<" and the tag end character ">", the matching starting position is assumed to exist outside the tag and the process proceeds to step 3804.

25 Step 3801:

This step checks whether or not the query term coincides with the character string from the matching starting position. In the case where the string

TOP SECRET - 0944950

character of the query term includes a plurality of bytes, the character string is checked byte by byte. In the case where the string character of the query term coincides with the character string from the matching string position, the process proceeds to step 3802. Otherwise, the process proceeds to step 3803.

Step 3802:

In the case where the query term is coincident with the string character from the matching starting position at step 3801, "matching" is assumed and the process is terminated.

Also, the end position of the matching character string is determined. The end position is assumed to be the number of bytes equal to the sum of the matching start character position and the character string length of the query term. The end position thus determined is used at step 3711.

Step 3803:

In the case where the query term is not coincident at step 3801, "no matching" is assumed and the process is terminated.

Step 3804:

In the case where the matching starting position exists outside the tag at step 3800, the process of searching outside the tag is performed. The out-of-tag search process will be explained with reference to Fig. 39.

Step 3805:

This step checks whether there exists in the HTML document a point matching with the query term at step 3804. In the case where there is any such a point,
5 the process proceeds to step 3807. Otherwise, the process proceeds to step 3806.

Step 3806:

In the case where the query term fails to match at step 3805, the process is terminated.

10 Step 3807:

In the case where the query term matches at step 3805, on the other hand, "matching" is assumed and the process is terminated.

Also, the end position of the matching
15 character string is determined. The end position is assumed to be the sum of the matching start character position and the position detected at step 3804 where the last matching character is described. The end position thus determined is used at step 3711.

20 The search inside the tag and search outside the tag were described above.

Now, the process of search outside the tag of step 3804 will be explained with reference to Fig. 39.

Step 3900:

25 This step checks whether or not there exists a query term in the HTML document. It is checked whether or not the character string of the query term coincides with the character string existing in the HTML document.

TOP SECRET 5-449550

5

Specifically, this process will be explained with reference to Fig. 34.

10

In the case under consideration, the query term is checked with the HTML document character by character.

20

25

5

10

15

Step 3902:

20

25

In the case where the character of the HTML document is the tag start character "<" at step 3902, the contents of the tag are skipped and the process

5

10

15

20

25

At step 3710, the highlight position information stored in the highlight position information storage area 2709 is read out.

The HTML document temporary storage area 2711 is secured for storing the HTML document with the highlight tag inserted therein.

10 The number of highlight tags is read out of
the highlight tag number storage area 2708. Also, the
character string length of the highlight start and end
tags is determined by detecting the tags from the
highlight tag number 3003 of the highlight position
15 information storage area 2709 and the highlight tag
character storage area 2710. Step 4002:

20 Step 4003:

Step 4004:

The data from i_cnt indicating the processed

position to the matching starting position are stored in the HTML temporary storage area 2711.

Specifically, in the case where the query term is a "feature article" in the HTML document 3400 of Fig.

5 34, the data from the head of the HTML document to the character "this month" before the characters "feature article" 3403 are all stored in the HTML document temporary storage area 2711.

Step 4005:

10 The highlight tag is stored in the matching strings position. The process for inserting the highlight tag will be explained later with reference to Fig. 41.

Step 4006:

15 The number of bytes from the head of the position at which the highlight end tag is inserted is substituted into i_cnt indicating the end position of the HTML document processing, and the process returns to step 4003.

20 Step 4007:

The data from i_cnt indicating the processed position of the HTML document to the end of the HTML document are stored in the HTML document temporary storage area 2711 and the process is terminated.

25 Now, the process of inserting the highlight tag processed at step 4005 will be explained with reference to Fig. 41.

In the case under consideration, it is checked

095447-02304
T08260"3449550

whether the matching strings position is outside or inside the tag, and a highlight tag is inserted before and after the matching strings position.

Step 4100:

5 It is checked whether or not the matching strings position of the HTML document is inside or outside the HTML tag.

 This check is carried out by a method in which the correspondence is taken between the start character
10 "<" and the end character ">" of the HTML tag in the HTML document up to the start of the matching string position, and it is determined whether the start of the matching strings position is inside or outside the tag. In the case where the start of the matching starting
15 position is located between the tag start character "<" and the tag end character ">", the start of the matching start position is assumed to exist inside the tag, and the process proceeds to step 4101. Otherwise, the start of the matching strings position is assumed to exist
20 outside the tag and the process proceeds to step 4110.

Step 4101:

 The character following the tag start character "<" is extracted and thus the tag type is acquired.

25 In the case of the HTML document 3400 of Fig. 34, for example, assume that the query term is "hitachi". The term "hitachi" can be acquired at 3409 in the HTML document 3400. Once the type of the HTML

096443-09304
TOP SECRET 52449660

Step 4102:

10

15

20

25

This HTML tag is linked to a URL (uniform resource locator) specified in the start tag when a

5

10

20

Step 4104:

25

In the case under consideration, the characters up to the last character ">" of the end tag are skipped thereby to produce the skipped data.

Specifically, in the case where the query term is "hitachi" for the HTML document 3400 of Fig. 34, the HTML tag is "A" (3410). The end tag is "" (3412).

Thus, the data up to "" (3412) is
5 acquired.
Step 4105:

In the case where a character can be inserted between the start tag and the end tag, the data up to the last character ">" of the tag is acquired.

10 Specifically, assuming that the query term for the HTML document 3400 of Fig. 34 is "imagefile.gif" (3411), the type of the HTML tag is "IMG" (3413) with the last tag character of ">" (3414). Therefore, the characters up to ">" (3414) are skipped, so that the
15 skipped data are obtained.
Step 4106:

The data acquired at steps 4104 and 4105 are inserted in the HTML document temporary storage area 2711. Also, the end position of the matching character
20 is determined. The end position is represented by the number of bytes of the position skipped at step 4104 or 4105.

Step 4107:

The start tag for highlight is inserted in the
25 HTML document temporary storage area 2711. The tag name written on the start tag 1 (3203) of Fig. 32 is inserted in the HTML document temporary storage area 2711.

In the case where the data stored in the

TOP SECRET 5644950

5 storage area. It is therefore possible to acquire
" <BLINK>" from (2) of Fig. 32. As a result, " <BLINK>" is
inserted in the HTML document temporary storage area
2711.

10 The matching character is inserted in the HTML
document temporary storage area 2711. Since the match-
ing character is rewritten, however, a rewrite mark is
stored. In the case under consideration, the data
stored in the rewrite mark storage area 3602 is
15 inserted.

Step 4109:

25 Step 4110:

In the case where the matching strings position is located outside the tag, the highlight tag

for out-of-tag application is inserted. This process will be described later with reference to Fig. 42.

Now, an explanation will be given of the process at step 4110 for inserting the highlight tag for out-of-tag application with reference to Fig. 42.

Step 4200:

For the tag "<A~>", the emphasized display is not reflected if the highlight tag is inserted at a point sandwiched between the start tag "<A~>" and the end tag "". In the case under consideration, it is checked whether or not the particular portion is surrounded by the start tag and the end tag, and it is determined at step 4201 whether or not the highlight tag can be inserted between the tags. The tags "HTML" and "<BODY>" exist for the HTML document and are used as tags surrounding the whole document, and therefore are not checked for the present purpose.

It is checked whether or not the position stored in the highlight-position-from-head information 3001 of Fig. 30 is surrounded by the start tag and the end tag of the HTML document. In the case where it is so surrounded, the process proceeds to step 4201. Otherwise, the process proceeds to step 4209.

Specifically, in the case where the query term is "HITACHI" for the HTML document 3400 of Fig. 34, "HITACHI" (3415) of the HTML document 3400 is extracted. Since this query term is surrounded by the tag "<A~>"

(3410) and "" (3412), the process proceeds to step 4201.

Step 4201:

It is checked whether or not a highlight tag
5 can be inserted before and after the matching character string.

In a checking method, the HTML tag surrounded
by the matching character string is extracted, and the
type of the extracted tag is compared with the tag
10 stored in the highlight tag no-insertion tag name
storage area 3603. In the case of coincidence, the
process proceeds to step 4202, and otherwise, to step
4209.

The highlight tag in the name written in the
15 highlight tag no-insertion tag name storage area 3603
cannot be inserted between the start tag and the end
tag.

The highlight tag is thus compared with the
HTML tag stored in the highlight tag no-insertion tag
20 name storage area 3603, and if coincident, the process
proceeds to step 4202. Otherwise, the process proceeds
to step 4209. The highlight tag no-insertion tag name
storage area 3603 is produced before data control (2705)
using the user interface.

25 Step 4202:

In the case where the highlight tag cannot be
used, the HTML document up to the last character ">" of
the end tag is skipped.

09964435-092804
T08260"56449660

In Fig. 34, in the case where the query term is "HITACHI", the HTML document up to "" (3412) is skipped.

Step 4203:

5 The HTML document skipped at step 4202 is stored in the HTML document temporary storage area 2711.

 In Fig. 34, in the case where the query term is "HITACHI", the data of "" (3412) is stored in the HTML document temporary storage area 2711 from the
10 (i_cnt)th character of the HTML document set at step 4002 or from the (i_cnt)th character of the HTML document updated at step 4006.

Step 4204:

 The start tag for highlight is inserted in the
15 HTML document temporary storage area 2711. In the case where the highlight position information storage area is 3402 and the highlight tag character storage area is located in (2) of Fig. 32, "<BLINK>" is extracted. Thus, in the case under consideration, "<BLINK>" is
20 inserted.

Step 4205:

 A redisplay mark is stored. Like step 4108, the HTML document stored in the rewrite mark storage area 3602 is read out, and stored in the HTML document
25 temporary storage area 2711.

Step 4206:

 The matching character string is inserted again in the HTML document temporary storage area 2711.

05564475-092804

Step 4207:

Step 4208:

In the case where a tag exists in the matching character string and all the matching characters are not yet stored, the process returns to step 4200. In the case where all the characters are stored, on the other

Step 4209:

Specifically, in the case where the query term for the HTML document 3400 of Fig. 34 is "feature article", the HTML document up to "this month" existing before the matching "feature" (3403) is inserted in the HTML document temporary storage area 2711.

As at step 4205, the start tag for highlight is stored in the HTML document temporary storage area

2711. In the case under consideration, "<BLINK>" is inserted.

Step 4211:

5 The matching character string is inserted in the HTML temporary storage area 2711.

In the case where a tag exists midway of the matching character string, however, the character string up to the point of the tag is inserted.

10 In the case where the query term is "feature article" for the HTML document 3400, for example, "</H1>" (3417) exists between "feature" (3403) and "article" (3416). In this case, therefore, "feature" is stored.

Step 4212:

15 The end tag for highlight is inserted in the HTML document temporary storage area 2711. In the case under consideration, "</BLINK>" is inserted.

Step 4213:

20 It is checked whether or not all the string characters of the query term are inserted in the HTML document temporary storage area 2711. Assume that the query term matches the character string of the HTML document when the HTML tag is removed, and that the HTML tag exists between the head of the matching position and
25 the character string having the length of the query term. Then, the HTML document up to the HTML tag is inserted in the HTML document temporary storage area

09064475-09001
1003250-5449550

2711 at step 4211. In this case, it is necessary to process the remaining matching characters from the HTML tag.

In the case where all the query terms are
5 inserted in the HTML document temporary storage area, the process is terminated. Also, in the case where the remaining matching characters from the HTML tag are processed, the process is returned to step 4200.

In the case where the query term is "feature
10 article" for the HTML document 3400 of Fig. 34, "</H1>" (3417) exists between "feature" (3403) and "article" (3416). Since "article" is not inserted but "feature" at step 4206, the process is returned to step 4200.

The above-mentioned process makes it possible
15 to insert a highlight tag in the HTML document matching with the query term and to display the highlight matching strings position on the web browser 2703 using the query term set by the client 2701.

An explanation was given that in the present
20 case, the HTML document is checked for a single query term and in the presence of a query term character string in the HTML document, the result of search is displayed on the web browser of the client 2701. Nevertheless, it is possible to search a plurality of
25 HTML documents for a single query term, to store the highlight position information equivalent to the number of the matching HTML documents and to produce a

09564475-092804
T08260-52449660

plurality of HTML documents collectively with the highlight tag stored therein.

It is also possible to search a plurality of HTML documents for a plurality of query terms, to store
5 the highlight position information corresponding to the number of the matching HTML documents and to produce a plurality of HTML documents collectively with the highlight tag stored therein.

Now, a sixth embodiment of the invention will
10 be explained.

The difference of this embodiment from the second embodiment is that a highlighting method as well as a query term can be defined in a query in the case where the query is matched. As a result, a highlighting
15 method can be specified for each arbitrary query.

The system configuration of this embodiment is identical to that of Fig. 1, except that the method of writing the query 103 is different. An example of the method of writing the query 103 according to this
20 embodiment will be described with reference to Fig. 43.

Fig. 43 shows an example query according to this embodiment. As shown in Fig. 43, a highlighting method like "{underline}" is specified after each query term or each element. The query in the second
25 embodiment is "specify element to be searched: query expression". The query according to the invention, on the other hand, is "element to be searched {highlighting method}: query expression with highlighting method".

It is possible to eliminate the specification of the highlighting method. When the specification of the highlighting method is eliminated, the highlighted display is carried out by the method shown in the second embodiment. Specifically, with regard to the portion for which the highlighting method is not written in the query, the definition of the highlighting method shown in Fig. 18 is read out, and the highlighted display carried out using the highlighting method described in the definition information.

Fig. 44 shows the contents stored as the matching information 4401 according to this embodiment. The difference from the matching information shown in Fig. 17 of the second embodiment lies in that the highlighting method 4403 as well as the matching condition 4402 is stored for each matching strings position. This information can be acquired by analyzing the above-described query with reference to Fig. 43 and reading out the information on the highlighting method written in the query.

Fig. 45 shows a method of generating a DTD for highlighted display according to this embodiment. In this embodiment, in view of the fact that the highlighting method may be altered each time of search, only the required element is added to generate a DTD for highlighted display for each highlighted display. In this case, not the query but the highlighting method is written directly in the DTD.

As shown in Fig. 45, in addition to the original DTD 1901 used for registration, a containing element for highlight has generated therein a DTD 4501 for highlighted display with the definition altered and
5 added in such a manner as to permit hierarchical specification of a subelement for highlight.

A method of producing a DTD will be described. First, in the case where the highlighting method 4403 is not described in the matching strings position information of Fig. 44, a highlighting method corresponding to
10 the matching condition is acquired from the definition of the highlighting method shown in Fig. 18. The element information is altered (4502) to make it possible to produce a content model of a highlighting method
15 occurring in the subelement of each element of the original DTD. Further, a hierarchical relation of the highlighting element that occurs is acquired from the hierarchical relation of the matching strings positions in the matching strings position information 4401.
20 Based on the hierarchical relation for highlighted display thus acquired, each highlighting element is rendered to have a subelement in the form of a highlighting subelement and a character string as a content model. In the absence of a highlighting
25 subelement, only a character string is caused to occur as a content model (4503).

The highlighting process according to this embodiment is not to make a highlighting element of a

108260-544355

query, but to generate a structured document for highlighted display describing a highlighting method and to generate a DTD for highlighted display. For this purpose, a structured document for display as shown in
5 Fig. 46 is produced according to this embodiment. This structured document for highlighted display is displayed in highlight as shown in Fig. 47.

According to this invention, when displaying the contents of the matching document as the result of
10 searching the structured document, it is possible to output a structured document with highlight information added thereto at a position matching the query term for each element. The highlighted display is made possible for any browser with the highlight information embedded
15 in the structured document but not dependent on the browser.

Different highlighting processes are possible according to such conditions as the importance and frequency of occurrence of each query or query term. As
20 to an crucial query term, therefore, a highlighting process specifying a high degree of weighting can be performed. Further, the description of a highlighting method in the query makes possible an arbitrary highlighted display for each user.

25 Furthermore, it is possible to extract only a subelement and output a structured document with the highlight information added thereto.

Also, a matching is easily secured in the case

where a document having a HTML tag indicating the document structure therein is searched for a character string, in the case where a character string coincident with a set query term exists in the HTML tag or in the
5 case where a query term is described over a HTML tag.
In addition, a matching character string can be displayed in highlight.

0994475-092801
TOP SECRET